

What is claimed is:

- 1 1. A system for analyzing unstructured documents for conceptual
2 relationships, comprising:
 - 3 a histogram module determining a frequency of occurrences of concepts in
4 a set of unstructured documents, each concept representing an element occurring
5 in one or more of the unstructured documents;
 - 6 a selection module selecting a subset of concepts out of the frequency of
7 occurrences, grouping one or more concepts from the concepts subset, and
8 assigning weights to one or more clusters of concepts for each group of concepts;
9 and
 - 10 a best fit module calculating a best fit approximation for each document
11 indexed by each such group of concepts between the frequency of occurrences
12 and the weighted cluster for each such concept grouped into the group of
13 concepts.
- 1 2. A system according to Claim 1, further comprising:
 - 2 an extraction module extracting features from each of the unstructured
3 documents and normalizing the extracted features into the concepts.
- 1 3. A system according to Claim 2, further comprising:
 - 2 a structured database storing the extracted features as uniquely identified
3 records.
- 1 4. A system according to Claim 1, further comprising:
 - 2 a visualization module visualizing the frequency of occurrences,
3 comprising at least one of creating a histogram mapping the frequency of
4 occurrences for each document in the unstructured documents set and creating a
5 corpus graph mapping the frequency of occurrence for all such documents in the
6 unstructured documents set.
- 1 5. A system according to Claim 1, further comprising:

2 a threshold comprising a median and edge conditions, each such concept
3 in the concepts subset occurring within the edge conditions.

1 6. A system according to Claim 1, further comprising:
2 an inner product module determining, for each group of concepts, the best
3 fit approximation as the inner product between the frequency of occurrences and
4 the weighted cluster for each such concept in the group of concepts.

1 7. A system according to Claim 6, wherein the inner product $d_{cluster}$ is
2 calculated according to the equation comprising:

3
$$d_{cluster} = \sum_{i=1}^n doc_{term_i} \cdot cluster_{term_i}$$

4 where $doc_{concept}$ represents the frequency of occurrence for a given concept in the
5 document and $cluster_{concept}$ represents the weight for a given cluster.

1 8. A system according to Claim 1, further comprising:
2 a control module iteratively re-determining the best fit approximation
3 responsive to a change in the set of unstructured documents.

1 9. A method for analyzing unstructured documents for conceptual
2 relationships, comprising:
3 determining a frequency of occurrences of concepts in a set of
4 unstructured documents, each concept representing an element occurring in one or
5 more of the unstructured documents;
6 selecting a subset of concepts out of the frequency of occurrences;
7 grouping one or more concepts from the concepts subset;
8 assigning weights to one or more clusters of concepts for each group of
9 concepts; and
10 calculating a best fit approximation for each document indexed by each
11 such group of concepts between the frequency of occurrences and the weighted
12 cluster for each such concept grouped into the group of concepts.

1 10. A method according to Claim 9, further comprising:

2 extracting features from each of the unstructured documents; and
3 normalizing the extracted features into the concepts.

1 11. A method according to Claim 10, further comprising:
2 storing the extracted features as uniquely identified records in a structured
3 database.

1 12. A method according to Claim 9, further comprising:
2 visualizing the frequency of occurrences, comprising at least one of:
3 creating a histogram mapping the frequency of occurrences for
4 each document in the unstructured documents set; and
5 creating a corpus graph mapping the frequency of occurrence for
6 all such documents in the unstructured documents set.

1 13. A method according to Claim 9, further comprising:
2 defining a threshold comprising a median and edge conditions, each such
3 concept in the concepts subset occurring within the edge conditions.

1 14. A method according to Claim 9, further comprising:
2 for each group of concepts, determining the best fit approximation as the
3 inner product between the frequency of occurrences and the weighted cluster for
4 each such concept in the group of concepts.

1 15. A method according to Claim 14, wherein the inner product $d_{cluster}$
2 is calculated according to the equation comprising:

$$d_{cluster} = \sum_{i=1}^n doc_{term_i} \cdot cluster_{term_i}$$

4 where $doc_{concept}$ represents the frequency of occurrence for a given concept in the
5 document and $cluster_{concept}$ represents the weight for a given cluster.

1 16. A method according to Claim 9, further comprising:
2 iteratively re-determining the best fit approximation responsive to a
3 change in the set of unstructured documents.

1 17. A computer-readable storage medium holding code for performing
2 the method according to Claims 9, 10, 11, 12, 13, 14, 15, or 16.

1 18. A system for dynamically evaluating latent concepts in
2 unstructured documents, comprising:

3 an extraction module extracting a multiplicity of concepts from a set of
4 unstructured documents into a lexicon uniquely identifying each concept and a
5 frequency of occurrence;

6 a frequency mapping module creating a frequency of occurrence
7 representation for each documents set, the representation providing an ordered
8 corpus of the frequencies of occurrence of each concept;

9 a concept selection module selecting a subset of concepts from the
10 frequency of occurrence representation filtered against a minimal set of concepts
11 each referenced in at least two documents with no document in the corpus being
12 unreferenced;

13 a group generation module generating a group of weighted clusters of
14 concepts selected from the concepts subset; and

15 a best fit module determining a matrix of best fit approximations for each
16 document weighted against each group of weighted clusters of concepts.

1 19. A system according to Claim 18, further comprising:

2 a histogram module creating a histogram mapping the frequency of
3 occurrence representation for each document in the documents set.

1 20. A system according to Claim 19, further comprising:

2 a data mining module mining the multiplicity of concepts from each
3 document as at least one of a noun, noun phrase and tri-gram.

1 21. A system according to Claim 19, further comprising:

2 a normalizing module normalizing the multiplicity of concepts into a
3 substantially uniform lexicon.

1 22. A system according to Claim 21, wherein the substantially uniform
2 lexicon is in third normal form.

1 23. A system according to Claim 18, further comprising:
2 a corpus mapping module creating a corpus graph mapping the frequency
3 of occurrence representation for all documents in the documents set.

1 24. A system according to Claim 18, further comprising:
2 a threshold module defining the pre-defined threshold as a median value
3 and a set of edge conditions and choosing those concepts falling within the edge
4 conditions as the concepts subset.

1 25. A system according to Claim 18, further comprising:
2 a cluster module naming one or more of the concepts within the concepts
3 subset to a cluster and assigning a weight to each concept with each such cluster.

1 26. A system according to Claim 25, further comprising:
2 a group module grouping one or more of the clusters into each such group
3 of weighted clusters of concepts.

1 27. A system according to Claim 18, further comprising:
2 a Euclidean module calculating a Euclidean distance between the
3 frequency of occurrence for each document and a corresponding weighted cluster.

1 28. A system according to Claim 18, further comprising:
2 a iteration module removing select documents from the documents set and
3 iteratively reevaluating the matrix of best fit approximations based on a revised
4 frequency of occurrence representation and concepts subset.

1 29. A system according to Claim 18, further comprising:
2 a structured database storing the lexicon, the lexicon comprising a
3 plurality of records each uniquely identifying one such concept and an associated
4 frequency of occurrence.

1 30. A system according to Claim 29, wherein the structured database is
2 an SQL database.

1 31. A method for dynamically evaluating latent concepts in
2 unstructured documents, comprising:
3 extracting a multiplicity of concepts from a set of unstructured documents
4 into a lexicon uniquely identifying each concept and a frequency of occurrence;
5 creating a frequency of occurrence representation for each documents set,
6 the representation providing an ordered corpus of the frequencies of occurrence of
7 each concept;
8 selecting a subset of concepts from the frequency of occurrence
9 representation filtered against a minimal set of concepts each referenced in at least
10 two documents with no document in the corpus being unreferenced;
11 generating a group of weighted clusters of concepts selected from the
12 concepts subset; and
13 determining a matrix of best fit approximations for each document
14 weighted against each group of weighted clusters of concepts.

1 32. A method according to Claim 31, further comprising:
2 creating a histogram mapping the frequency of occurrence representation
3 for each document in the documents set.

1 33. A method according to Claim 32, further comprising:
2 mining the multiplicity of concepts from each document as at least one of
3 a noun, noun phrase and tri-gram.

1 34. A method according to Claim 32, further comprising:
2 normalizing the multiplicity of concepts into a substantially uniform
3 lexicon.

1 35. A method according to Claim 34, wherein the substantially
2 uniform lexicon is in third normal form.

1 36. A method according to Claim 31, further comprising:
2 creating a corpus graph mapping the frequency of occurrence
3 representation for all documents in the documents set.

1 37. A method according to Claim 31, further comprising:
2 defining the pre-defined threshold as a median value and a set of edge
3 conditions; and
4 choosing those concepts falling within the edge conditions as the concepts
5 subset.

1 38. A method according to Claim 31, further comprising:
2 naming one or more of the concepts within the concepts subset to a
3 cluster; and
4 assigning a weight to each concept with each such cluster.

1 39. A method according to Claim 38, further comprising:
2 grouping one or more of the clusters into each such group of weighted
3 clusters of concepts.

1 40. A method according to Claim 31, further comprising:
2 calculating a Euclidean distance between the frequency of occurrence for
3 each document and a corresponding weighted cluster.

1 41. A method according to Claim 31, further comprising:
2 removing select documents from the documents set; and
3 iteratively reevaluating the matrix of best fit approximations based on a
4 revised frequency of occurrence representation and concepts subset.

1 42. A method according to Claim 31, further comprising:
2 storing the lexicon in a structured database, the lexicon comprising a
3 plurality of records each uniquely identifying one such concept and an associated
4 frequency of occurrence.

1 43. A method according to Claim 42, wherein the structured database
2 is an SQL database.

1 44. A computer-readable storage medium holding code for performing
2 the method according to Claims 31, 32, 33, 34, 36, 37, 38, 39, 40, 41, or 42.